# Optimization & Algorithms

# MEEC

---

**Project Report**

---

**Group: 4**

Alexandre Leal (103098)                    alexandre.b.leal@tecnico.ulisboa.pt

Diogo Ramos (100299)                    diogo.ramos@tecnico.ulisboa.pt

Diogo Sampaio (103068)                    diogo.sampaio@tecnico.ulisboa.pt

Francisco Tavares (103402)        francisco.carreira.tavares@tecnico.ulisboa.pt

**2024/2025 – 1º Semestre, P1**

# 1  Task 1: Theoretical Task

In this task, we demonstrate that the function

$$f_D(w_0, w) = \frac{1}{N} \sum_{n=1}^{N} 1_{R_-} \left( y_n C_{w_0,w}(x_n) \right)$$

is not convex in the simplified case where $N = 1$ and $D = 1$. For $N = 1$, it simplifies to:

$$f_D(w_0, w) = 1_{R_-} \left( y_n C_{w_0,w}(x_n) \right).$$

We define two functions:

$$g(u) = 1_{R_-}(u), \quad f_2(w_0, w) = y_n C_{w_0,w}(x_n).$$

Thus, we express $f_D$ as the composition:

$$f_D(w_0, w) = g \circ f_2.$$

If $g(u)$ is not convex, then $f_D$ is also not convex if $f_2$ maps to a region where $g(u)$ is not convex. Observing figure 1, we can see that $g(u)$ is not convex.

To demonstrate, using the definition of convexity consider $x = -1$, $y = 1$, and $\alpha = 0.25$:

$$
\begin{aligned}
g((1-\alpha)x + \alpha y) &\leq (1-\alpha)g(x) + \alpha g(y) \\
&\iff \quad g(0.75 \cdot (-1) + 0.25 \cdot 1) \leq 0.75 g(-1) + 0.25 g(1) \\
&\iff \quad g(-0.5) \leq 0.75 \cdot 1 + 0.25 \cdot 0 \\
&\iff \quad 1 \leq 0.75.
\end{aligned}
$$

This is false; hence, $g(u)$ is not convex, and so $f_D$ is also not convex.

**Note:** The function $f_2$ maps $\mathbb{R}$ to $\{+1, -1\}$ since $y_n$ and $C_{w_0,w}$ belong to $\{+1, -1\}$, validating the points used in the example.
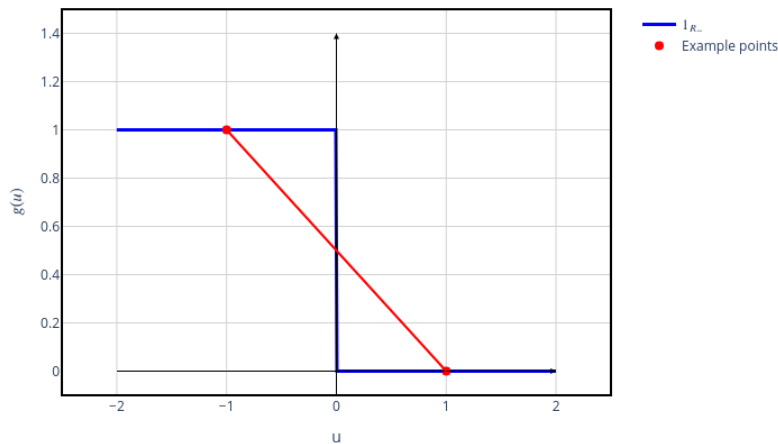


**Figure 1:** Plot of the function $g(u)$

# 2   Task 2. [Theoretical Task]

In this task we were asked to show that the function $1_{R_-}$ is majorized by $h$, and that $1_{R_-} \le h(u)$ holds for all $u \in \mathbb{R}$. Furthermore we showed that $h$ is a convex function.
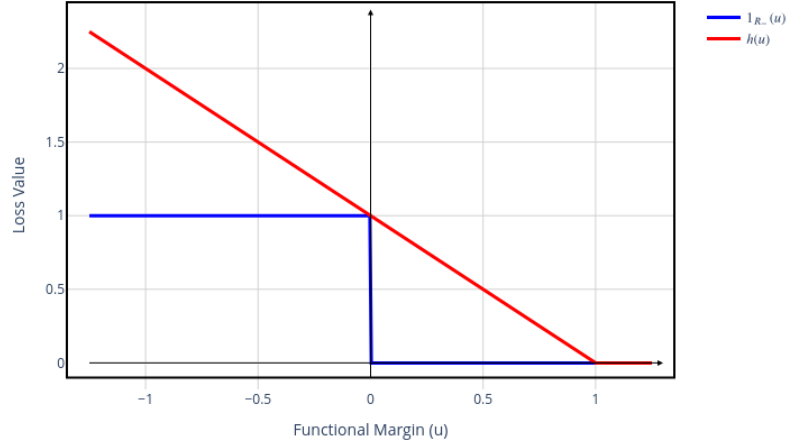


**Figure 2:** Plot of the indicator function $1_{R_-}(u)$ & hinge loss function $h(u)$

We want to prove that $1_{R_-}(u) \le h(u)$ for all $u \in \mathbb{R}$. We can describe $1_{R_-}(u)$, and $h(u)$ as:

$$h(u) = \begin{cases} 0, & u \le 1 \\ 1 - u, & u > 1 \end{cases}$$

$$1_{R_-}(u) = \begin{cases} 1, & u < 0 \\ 0, & u \ge 0 \end{cases}$$

For $u < 0$:

$$1_{R_-}(u) \le h(u) \Longleftrightarrow 1 \le 1 - u$$
$$\Longleftrightarrow 0 \le -u \Longleftrightarrow u \le 0 \quad True$$

For $u \in [0; 1[$:

$$1_{R_-}(u) \le h(u) \Longleftrightarrow 0 \le 1 - u$$
$$\Longleftrightarrow u \le 1 \quad True$$

For $u \ge 1$:

$$1_{R_-}(u) \le h(u) \Longleftrightarrow 0 \le 0 \quad True$$

To show that $h$ is a convex function we can observe figure 3, it´s easy to see that no matter which 2 points in $h$ we choose, if we draw a straight line between them, $h$ will be always lower or equal then any point in that line, therefore, $h$ is convex.
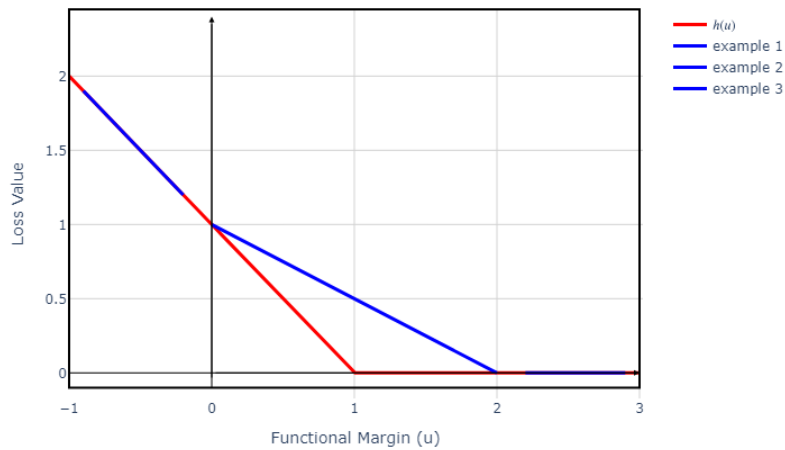


**Figure 3:** Plot of the hinge loss function $h(u)$ and some examples to show convexity

# 3  Task 3: Theoretical Task

We are asked to show that the function

$$g_D(\omega_0, \omega) = \frac{1}{N} \sum_{n=1}^{N} h(y_n(\omega_0 + x_n^T \omega))$$

is convex for any $N$ and $D$.

The function $g_D(\omega_0, \omega)$ can be written as a sum of convex functions in the form:

$$g_D(\omega_0, \omega) = \frac{1}{N} \sum_{n=1}^{N} h(y_n(\omega_0 + x_n^T \omega)) = \frac{1}{N} h(y_1(\omega_0 + x_1^T \omega)) + \cdots + \frac{1}{N} h(y_N(\omega_0 + x_N^T \omega)).$$

If $f_1, \ldots, f_N$ are convex, then $g_D(\omega_0, \omega)$ is convex because a sum of convex functions is convex [slide 12 of module 3 of the theoretical slides]. Since all $f_n$ have the same form, we only need to prove that $f_n$ is convex for all $n \in \{1, \ldots, N\}$.

We define

$$f_n = h(y_n(\omega_0 + x_n^T \omega)),$$

which can be written as $f_n = h(g(\omega_0, \omega))$. To show $f_n$ is convex, we need to check two conditions: 1. $h(u)$ is convex, and 2. $g(\omega_0, \omega)$ is affine [slide 12 of module 3 of the theoretical slides].

First, $g(\omega_0, \omega)$ can be written as

$$g(\omega_0, \omega) = y_n(\omega_0 + x_n^T \omega) = \begin{bmatrix} y_n & y_n x_n^T \end{bmatrix} \begin{bmatrix} \omega_0 \\ \omega \end{bmatrix},$$

where

$$A = \begin{bmatrix} y_n & y_n x_n^T \end{bmatrix} \in \mathbb{R}^{1 \times (D+1)}, \quad b = 0, \quad \text{and} \quad x = \begin{bmatrix} \omega_0 \\ \omega \end{bmatrix} \in \mathbb{R}^{D+1}.$$

Thus, $g(\omega_0, \omega)$ is an affine map of the form $g(\omega_0, \omega) = Ax + b$, where $A$ is a matrix and $b$ is a vector.

Since affine functions preserve convexity, and we are given that $h(u)$ is convex (as shown in Task 2), it follows that $f_n = h(g(\omega_0, \omega))$ is convex. Therefore, $g_D(\omega_0, \omega)$ is convex.

# 4   Task 4. [Theoretical Task]

We can write $g$ as $g_d(w_0, w) + p \cdot (\|w\|_2)^2$, where $p > 0$, which is in the form $t_1 f_1 + t_2 f_2$. Therefore, if $g_d$ is convex (proven in Task 3) and $(\|w\|_2)^2$ is convex, then $g$ is convex.

We can write $(\|w\|_2)^2$ as $f \circ h(w)$, with $h(w) = \|w\|_2$ and $f(u) = (u_+)^2$. We can say this because for $u \geq 0$, $u_+ = u$ and we know that $h(w) \geq 0$ for all $w$. $f \circ h(w)$ is convex if $f$ is convex and non-decreasing and $h$ is convex. $h$ is an elementary convex function (the $l_p$ norm) [slide 7 of module 3 of the theoretical slides], so it is convex.

$f(u)$ can be described as:

$$
\begin{cases}
f(u) = u^2, u \geq 0 \\
f(u) = 0, u < 0
\end{cases}
\tag{1}
$$

Because $f$ is continuous, if both segments are non-decreasing, then $f$ is also non-decreasing. For $u < 0$, $f$ is constant, so it is non-decreasing. For $u \geq 0$, $f'(u) = 2u$. Since $u \geq 0$, it means the slope of $f$ is always positive, which implies that $f$ is non-decreasing. This means $(\|w\|_2)^2$ is convex and therefore $g$ is also convex.

# 5 Task 5. [Theoretical Task]

To determine if the function $g$ is strongly convex we will evaluate the case where $w = 0$ and $N = 1$. In this range of values we get $g(\omega_0, \omega) = h(y\omega_0)$. If we analyze the range where $yw_0 \geq 1$, g is a null constant, which means it is not strictly nor strongly convex. Equivalently, if we consider $f = g - \frac{-m||\omega_0||^2}{2} = \frac{-m||\omega_0||^2}{2}$, in the same range, for all positive values of m, f is clearly not convex, which equates to g not being strongly convex.

Now using the definition of strongly convex, for $g$ to be SCVX,

$$g((1 - \alpha)x + \alpha y) \leq (1 - \alpha)g(x) + \alpha g(y)\frac{-n\alpha(1 - \alpha)}{2}||x - y||^2 \tag{2}$$

Let's take for example $x = 20$ and $y = 30$ and $\alpha = 0.5$,

$$g(10 + 15) \leq 0.5g(20) + 0.5g(y)\frac{-n}{8}||x - y||^2 \tag{3}$$

Knowing that for $x > 1$, $g(x) = 0$, $x = 20, y = 30$

$$0 \leq -n\frac{||x - y||^2}{8} \qquad \frac{||x - y||^2}{8} > 0 \tag{4}$$

$0 \leq -n$ is false for all positive values of n therefore g is not strongly convex.

Furthermore, looking at the graph of the hinge loss function at the end of Task 2, we can also see the constant part of the function mentioned above, proving that g is not strongly convex.

# 6    Task 6. [Numerical Task]

In this task, we solve the following optimization problem numerically, using CVX in Python. The problem involves classifying images of handwritten digits from the MNIST dataset, where the digits 0 and 1 are considered. We use a linear classifier based on the hinge loss function with regularization, as stated in equation (5) of the problem.

The optimization problem we are solving is defined as follows:

$$\min_{w_0,w} \left[ \frac{1}{N} \sum_{n=1}^{N} \max(0, 1 - y_n(w^T x_n + w_0)) + \lambda \|w\|_2^2 \right] \tag{5}$$

The first term in the objective function represents the hinge loss, and the second term is an $\ell_2$ regularization term to prevent overfitting.

After solving the optimization problem, the following optimal parameters were obtained:

$$w_0 \approx 0.255$$

$$w \approx [0, 0, 0, \ldots, -1.76 \times 10^{-10}, -2.65 \times 10^{-10}, \ldots, 2.43 \times 10^{-3}, \ldots]$$

We also evaluated the classifier error rate $f_D$ on both the training and test datasets using the function:

$$f_D = \frac{1}{N} \sum_{n=1}^{N} \mathbf{1}(y_n \neq \text{sign}(w^T x_n + w_0))$$

The error rates were as follows for $\lambda$:

- Training dataset error rate: 0.00%

- Test dataset error rate: 0.12%

In addition to that, we also changed $\lambda$ to 0.5 to verify our code and noticed, as expected, that $f_d$ evaluates to 0.25% for the training set and 0.25% for the test set.

# 7   Task 7. [Theoretical Task]

## 7.1   Task 7 a.

This task is divided into two parts. First, we are asked to show that

$$\tilde{x} = x - P\operatorname{sgn}(yw)$$

solves the following equation:

$$
\begin{aligned}
&\underset{\tilde{x}}{\text{minimize}}\ y(w_0 + \tilde{x}^T w) \\
&\text{subject to } |\tilde{x}_d - x_d| \le P, \quad \text{for } 1 \le d \le D.
\end{aligned}
\tag{6}
$$

To solve this task, we started by expanding the cost function:

$$f(\tilde{x}) = y(\omega_0 + \tilde{x}^T\omega) = y\omega_0 + y\tilde{x}^T\omega,$$

Since $y\omega_0$ is a constant, we wish to minimize $y\sum_{d=1}^{D}\tilde{x}_d\omega$ where $\tilde{x}_d \in [x_d - P; x_d + P]$.

- If $y\omega_d > 0$, we should pick the smallest $\tilde{x}_d$, $\tilde{x}_d = x_d - P$.

- If $y\omega_d < 0$, we should pick the largest $\tilde{x}_d$, $\tilde{x}_d = x_d + P$.

- If $y\omega_d = 0$, any value of $\omega_d$ works because $y\tilde{x}_d\omega = 0$

$$
\tilde{x}_d = \begin{cases}
x_d - P, & \text{if } y\omega_d > 0 \\
x_d + P, & \text{if } y\omega_d < 0 \\
\text{any value}, & \text{if } y\omega_d = 0
\end{cases}
$$

From this analysis, we conclude that the optimal solution is given by $\tilde{x} = x - P\operatorname{sgn}(yw)$
Next, we need to verify whether this solution satisfies the specified constraints:

$$
\begin{aligned}
&|\tilde{x}_d - x_d| \le P, \quad d \in \{1, 2, 3\} \\
&|x - P\operatorname{sgn}(yw) - x| = |-P\operatorname{sgn}(yw)| = P \le P
\end{aligned}
\tag{7}
$$

This confirms that $\tilde{x} = x - P\operatorname{sgn}(yw)$ solves the optimization problem stated in (1).

## 7.2   Task 7 b.

In the second part of the task, we are asked to show that for $\tilde{x} = x - P \operatorname{sgn}(yw)$, the cost function in (1) evaluates to:

$$y(\omega_0 + x^T\omega) - P||y\omega||_1.$$

Substituting $\tilde{x}$ into the cost function:

$$\begin{aligned}
f(\tilde{x}) &= y\left(w_0 + (x - P\operatorname{sgn}(yw))^T w\right) \\
&\Longleftrightarrow y\left(w_0 + x^T w - P\operatorname{sgn}(yw)^T w\right) \\
&\Longleftrightarrow y\left(w_0 + x^T w - P\sum_{d=1}^{D} \operatorname{sgn}(yw_d)w_d\right) \\
&\Longleftrightarrow y(w_0 + x^T w) - P\sum_{d=1}^{D} |yw_d| \\
&\Longleftrightarrow y(w_0 + x^T w) - P||yw||_1.
\end{aligned}$$

which confirms the evaluation of the cost function.

Thus, we conclude that the derived expression correctly evaluates the cost function for $\tilde{x} = x - P\operatorname{sgn}(yw)$.

# 8    Task 8. [Numerical Task]

Using the result of Task 7, we replaced $(x_n, y_n)$ by $(\tilde{x}_n, y_n)$ with $\tilde{x}_n = x - P\operatorname{sgn}(yw)$ and evaluated the function $f_d$ defined in (2), obtaining the classifier of Task 6's error rate on the attacked test dataset. With $\lambda = 0.1$, as used in task 6, and P=0.18, $f_d$ evaluates 43.56% . We also changed $\lambda$ to 0.5 to verify our code and noticed, as expected, that $f_d$ evaluates 21.88%, which is very close to the expected value 21.9%

# 9   Task 9. [Numerical Task]

In this section we want to solve the problem of minimizing the effect of an attack on our classifier. This corresponds to solving the problem defined as:

$$\min_{w_0,w} \left[ \frac{1}{N} \sum_{n=1}^{N} \max(0, 1 - (y_n(w^T x_n + w_0) - P||y_n w||_1)) + \lambda \|w\|_2^2 \right] \quad (8)$$

with P = 0.18 and $\lambda = 0.1$.

On the training dataset $f_D$ evaluates to 0.75% and in the test dataset $f_D$ evaluates to 0.44%. As for the attacked dataset, $f_D$ evaluates to 2.19%.

We can observe that, while in the test and training sets, the classifier performs slightly worse, which is expected since it is not the ideal classifier for unaltered data, we see a large improvement on the performance on the attacked dataset, since this is the loss function we strive to minimize with this formulation.

# 10 Task 10. [Numerical Task]

In Task 10, we aim to solve the problem of fitting a piecewise-linear signal to a set of noisy measurements. The signal is modeled as a weighted combination of linear models, and we need to optimize the parameters of these models to minimize the fitting error.

The goal is to minimize the following objective function $f$:

$$\min_{s_1,r_1,\ldots,s_K,r_K,u_1,v_1,\ldots,u_{K-1},v_{K-1}} \sum_{n=1}^{N} \left(\hat{y}(x_n) - y_n\right)^2 \tag{9}$$

where $\hat{y}(x_n)$ is the predicted value at time $x_n$, given by:

$$\hat{y}(x_n) = \sum_{k=1}^{K} w_k(x_n)\hat{y}_k(x_n) \tag{10}$$

and $y_n$ represents the measured values.

Each linear model $\hat{y}_k(x)$ is defined as:

$$\hat{y}_k(x) = s_k x + r_k \tag{11}$$

where $s_k$ is the slope and $r_k$ is the intercept for the $k$-th model.

The weights $w_k(x)$, which determine how the models are combined, are given by the softmax function:

$$w_k(x) = \frac{e^{u_k x + v_k}}{\sum_{j=1}^{K} e^{u_j x + v_j}} \tag{12}$$

These weights $w_k(x)$ ensure that the combination of linear models is smooth across the different regions of the signal.

When solving for the gradient of the objective function, we first get:

$$\nabla f = 2 \sum_{n=1}^{N} \left(\hat{y}(x_n) - y_n\right) \cdot \nabla \hat{y}(x_n) \tag{13}$$

where $\nabla \hat{y}(x_n)$ represents the Jacobian matrix $J$.

The Jacobian matrix $J$ consists of the partial derivatives of $\hat{y}(x_n)$ with respect to the variables $s_k, r_k, u_k,$ and $v_k$:

$$J_{n,p} = \frac{\partial \hat{y}(x_n)}{\partial p} \tag{14}$$

where $p$ can represent $s_k, r_k, u_k,$ or $v_k$.

To compute the partial derivatives, we need to evaluate the derivatives of $w_k(x)$ and $w_k(j)$:

1. Derivative of $w_k(x)$ with respect to $u_k$:

$$\frac{\partial w_k(x)}{\partial u_k} = \frac{\partial}{\partial u_k}\left(\frac{e^{u_k x + v_k}}{\sum_{i=1}^{K} e^{u_i x + v_i}}\right)$$

$$= \frac{\frac{\partial}{\partial u_k}\left(e^{u_k x + v_k}\right) \cdot \left(\sum_{i=1}^{K} e^{u_i x + v_i}\right) - e^{u_k x + v_k} \cdot \frac{\partial}{\partial u_k}\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2}$$

$$= \frac{e^{u_k x + v_k} \cdot x \cdot \left(\sum_{i=1}^{K} e^{u_i x + v_i}\right) - e^{u_k x + v_k} \cdot \left(e^{u_k x + v_k} \cdot x\right)}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2}$$

$$= \frac{e^{u_k x + v_k} \cdot x \cdot \left(\sum_{i=1}^{K} e^{u_i x + v_i} - e^{u_k x + v_k}\right)}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2}$$

$$= w_k(x) \cdot x \cdot (1 - w_k(x)) \tag{15}$$

2. Derivative of $w_k(x)$ with respect to $v_k$:

$$\frac{\partial w_k(x)}{\partial v_k} = \frac{\partial}{\partial v_k}\left(\frac{e^{u_k x + v_k}}{\sum_{i=1}^{K} e^{u_i x + v_i}}\right)$$

$$= \frac{e^{u_k x + v_k} \cdot \left(\sum_{i=1}^{K} e^{u_i x + v_i}\right) - e^{u_k x + v_k} \cdot \frac{\partial}{\partial v_k}\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2}$$

$$= \frac{e^{u_k x + v_k} \cdot \left(\sum_{i=1}^{K} e^{u_i x + v_i}\right) - e^{u_k x + v_k} \cdot \left(e^{u_k x + v_k}\right)}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2}$$

$$= \frac{e^{u_k x + v_k} \cdot \left(\sum_{i=1}^{K} e^{u_i x + v_i} - e^{u_k x + v_k}\right)}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2}$$

$$= w_k(x) \cdot (1 - w_k(x)) \tag{16}$$

3. Derivative of $w_j(x)$ with respect to $u_k$ (for $j \neq k$):

$$\frac{\partial w_j(x)}{\partial u_k} = \frac{\partial}{\partial u_k}\left(\frac{e^{u_j x + v_j}}{\sum_{i=1}^{K} e^{u_i x + v_i}}\right)$$

$$= -\frac{e^{u_j x + v_j}}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2} \cdot e^{u_k x + v_k} \cdot x$$

$$= -w_j(x) \cdot w_k(x) \cdot x \tag{17}$$

4. Derivative of $w_j(x)$ with respect to $v_k$ (for $j \neq k$):

$$\frac{\partial w_j(x)}{\partial v_k} = \frac{\partial}{\partial v_k}\left(\frac{e^{u_j x + v_j}}{\sum_{i=1}^{K} e^{u_i x + v_i}}\right)$$

$$= -\frac{e^{u_j x + v_j}}{\left(\sum_{i=1}^{K} e^{u_i x + v_i}\right)^2} \cdot e^{u_k x + v_k}$$

$$= -w_j(x) \cdot w_k(x) \tag{18}$$

Using the results obtained from the derivatives above, we can calculate the following partial derivatives:

1. For $u_k$:

$$J_{n,k} = \frac{\partial \hat{y}(x_n)}{\partial u_k} = \frac{\partial}{\partial u_k}\left(\sum_{k=1}^{K} w_k(x_n)\hat{y}_k(x_n)\right)$$

$$= \hat{y}_k(x_n)\frac{\partial w_k(x_n)}{\partial u_k} + \sum_{j \neq k} \hat{y}_j(x_n)\frac{\partial w_j(x_n)}{\partial u_k}$$

$$= \hat{y}_k(x_n) \cdot w_k(x_n)(1 - w_k(x_n)) \cdot x_n - \sum_{j \neq k} \hat{y}_j(x_n) \cdot w_j(x_n)w_k(x_n) \cdot x_n$$

$$= w_k(x_n) \cdot x_n \cdot \left(\hat{y}_k(x_n)(1 - w_k(x_n)) - \sum_{j \neq k} w_j(x_n) \cdot \hat{y}_j(x_n)\right) \tag{19}$$

2. For $v_k$:

$$J_{n,(K-1)+k} = \frac{\partial \hat{y}(x_n)}{\partial v_k} = \frac{\partial}{\partial v_k}\left(\sum_{k=1}^{K} w_k(x_n)\hat{y}_k(x_n)\right)$$

$$= \hat{y}_k(x_n)\frac{\partial w_k(x_n)}{\partial v_k} + \sum_{j \neq k} \hat{y}_j(x_n)\frac{\partial w_j(x_n)}{\partial v_k}$$

$$= \hat{y}_k(x_n) \cdot w_k(x_n)(1 - w_k(x_n)) - \sum_{j \neq k} \hat{y}_j(x_n) \cdot w_j(x_n)w_k(x_n)$$

$$= w_k(x_n) \cdot \left(\hat{y}_k(x_n)(1 - w_k(x_n)) - \sum_{j \neq k} w_j(x_n) \cdot \hat{y}_j(x_n)\right) \tag{20}$$

3. For $s_k$:

$$J_{n,2(K-1)+k} = \frac{\partial \hat{y}(x_n)}{\partial s_k} = \frac{\partial}{\partial s_k}\left(w_k(x_n) \cdot (s_k x_n + r_k)\right)$$

$$= w_k(x_n) \cdot \frac{\partial}{\partial s_k}(s_k x_n + r_k)$$

$$= w_k(x_n) \cdot x_n \tag{21}$$

4. For $r_k$:

$$J_{n,(2(K-1)+K)+k} = \frac{\partial \hat{y}(x_n)}{\partial r_k} = \frac{\partial}{\partial r_k}\left(w_k(x_n) \cdot (s_k x_n + r_k)\right)$$

$$= w_k(x_n) \cdot \frac{\partial}{\partial r_k}(s_k x_n + r_k)$$

$$= w_k(x_n) \tag{22}$$

Finally, we can summarize the gradient as:

$$\nabla f = 2\sum_{n=1}^{N}\left(\hat{y}(x_n) - y_n\right)J \tag{23}$$

## 10.1  Results for Optimization Variables

After implementing the LM method and running the optimization, the algorithm converged to the following values for the variables:

$$u = \begin{bmatrix} -13.077 \\ -80.105 \\ 88.202 \end{bmatrix} \qquad v = \begin{bmatrix} -2.346 \\ -415.562 \\ -184.845 \end{bmatrix} \qquad s = \begin{bmatrix} -1.202 \\ 3.265 \\ -4.000 \\ 12.126 \end{bmatrix} \qquad r = \begin{bmatrix} -0.077 \\ 23.059 \\ 19.877 \\ -2.653 \end{bmatrix}$$
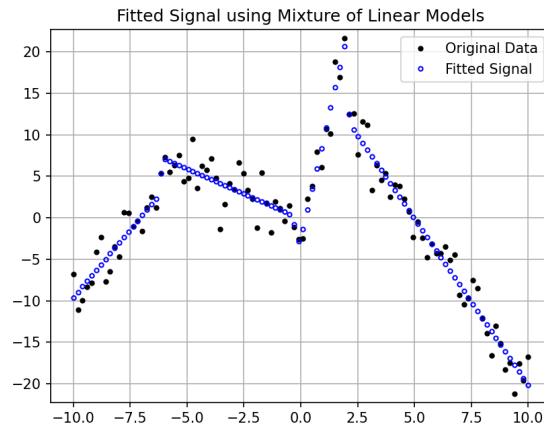
## 10.2  Plots



**Figure 4:** The fitted signal using the mixture of linear models. Each black dot represents a measurement, and the blue circles represent the output of the linear mixture.
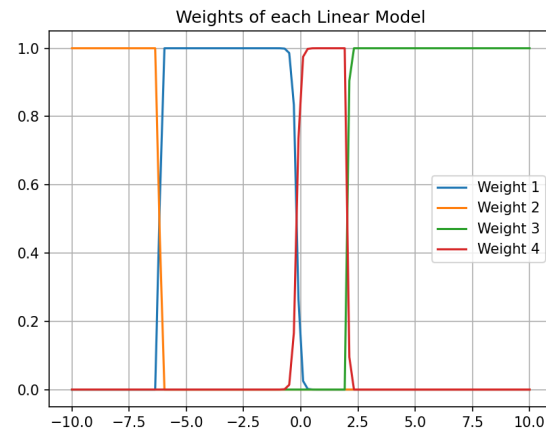
**Figure 5:** The weights $w_k(x)$ for each linear model. Each color represents a different weight function.
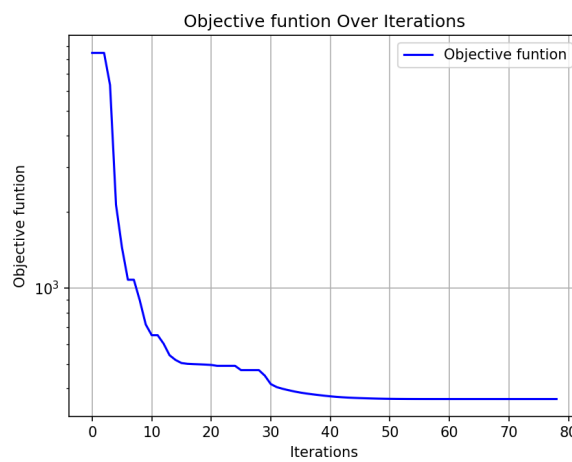


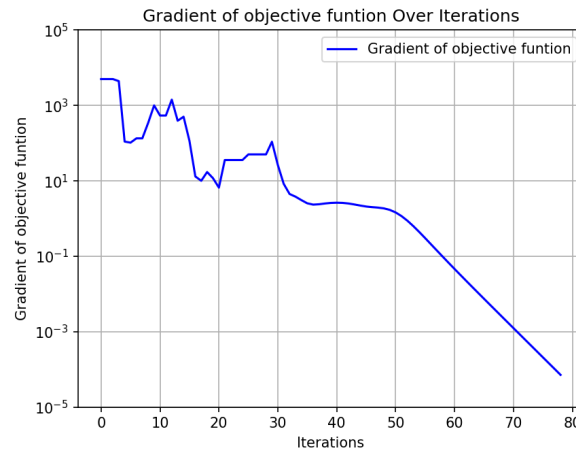**Figure 6:** Values of the objective function across the LM iterations.

**Figure 7:** The norm of the gradient of the objective function across the LM iterations.

## 10.3   Comment on the Results

The Levenberg-Marquardt method converged effectively after 78 iterations, successfully bringing the gradient norm below the set threshold of $\epsilon = 10^{-4}$. As a result, the fitted signal aligns closely with the noisy measurements we started with. The weight functions behaved as anticipated, with each model activating in different regions of the signal, contributing to the overall fit. Throughout the optimization process, the objective function consistently decreased, demonstrating strong convergence behavior, and we observed a clear reduction in the gradient norm over time.

In the end, the final fitted parameters provided an accurate representation of the signal, effectively capturing the distinct linear segments and integrating them into a coherent overall model.

# 11 Code Developed

## 11.1 Python Code For Task 6

```python
import cvxpy as cp
import numpy as np
import scipy.io
import matplotlib.pyplot as plt

# Open the data file .mat
data = scipy.io.loadmat('classifier_dataset.mat')

# Extract the variables from the data
traindataset = data['traindataset']
testdataset = data['testdataset']

trainlabels = data['trainlabels'].flatten()  # Flatten to convert into 1D
    array
testlabels = data['testlabels'].flatten()    # Flatten to convert into 1D
    array

# Save the number of rows in N (number of samples) and columns in D (number
    of features)
N, D = traindataset.shape

# Regularization parameter
ro = 0.1

# Define the optimization variables
w0 = cp.Variable()        # Bias term (scalar)
w = cp.Variable(D)        # Weights vector (D-dimensional)

# Define the hinge loss function h(u) = max(0, 1 - u)
hinge_losses = cp.pos(1 - cp.multiply(trainlabels, traindataset @ w + w0))

# Define the objective function: hinge loss + regularization term
objective = cp.Minimize((1/N) * cp.sum(hinge_losses) + ro * cp.norm(w, 2)
    **2)

# Define the problem and solve it
problem = cp.Problem(objective)
problem.solve()

# Retrieve the optimal parameters
w_optimal = w.value
w0_optimal = w0.value

# Function to evaluate the classifier error rate fD on a given dataset
def evaluate_error_rate(dataset, labels, w0_opt, w_opt):
    predictions = np.sign(dataset @ w_opt + w0_opt)  # Classify
    misclassifications = np.sum(predictions != labels)  # Count errors
    error_rate = misclassifications / len(labels)  # Calculate error rate
    return error_rate
```

```
47  # Evaluate the classifier error rate on the training dataset
48  train_error_rate = evaluate_error_rate(traindataset, trainlabels, w0_optimal
       , w_optimal)
49  print(f"Training dataset error rate: {train_error_rate * 100:.2f}%")
50
51  # Evaluate the classifier error rate on the test dataset
52  test_error_rate = evaluate_error_rate(testdataset, testlabels, w0_optimal,
       w_optimal)
53  print(f"Test dataset error rate: {test_error_rate * 100:.2f}%")
```

**Listing 1:** Python Code to Solve Task 6

## 11.2 Python Code For Task 8

```python
import cvxpy as cp
import numpy as np
import scipy.io
import matplotlib.pyplot as plt

# Open the data file .mat
data = scipy.io.loadmat('classifier_dataset.mat')

# Extract the variables from the data
traindataset = data['traindataset']
testdataset = data['testdataset']

trainlabels = data['trainlabels'].flatten()  # Flatten to convert into 1D
    array
testlabels = data['testlabels'].flatten()    # Flatten to convert into 1D
    array

# Save the number of rows in N (number of samples) and columns in D (number
    of features)
N, D = traindataset.shape

# Regularization parameter
ro = 0.1

# Define the optimization variables
w0 = cp.Variable()         # Bias term (scalar)
w = cp.Variable(D)         # Weights vector (D-dimensional)

# Define the hinge loss function h(u) = max(0, 1 - u)
hinge_losses = cp.pos(1 - cp.multiply(trainlabels, traindataset @ w + w0))

# Define the objective function: hinge loss + regularization term
objective = cp.Minimize((1/N) * cp.sum(hinge_losses) + ro * cp.norm(w, 2)
    **2)

# Define the problem and solve it
problem = cp.Problem(objective)
problem.solve()

# Retrieve the optimal parameters
w_optimal = w.value
w0_optimal = w0.value

# Function to evaluate the classifier error rate fD on a given dataset
def evaluate_error_rate(dataset, labels, w0_opt, w_opt):
    predictions = np.sign(dataset @ w_opt + w0_opt)  # Classify
    misclassifications = np.sum(predictions != labels)  # Count errors
    error_rate = misclassifications / len(labels)  # Calculate error rate
    return error_rate

P=0.18
```

```
50  ####################################################
51  #make the x˜ vector (atacker vector)
52  ####################################################
53  N2, D2 = testdataset.shape
54
55  x_attack_final=np.empty(shape=(N2, D2))
56
57  #For each sample, calculate x˜
58  for k in range(0,N2):
59      #calculate y*w
60      yw=testlabels[k]*w_optimal
61      i=0
62
63      #aply sign function to all elements in the vector yw
64      for num in yw:
65          if num>=0:
66              yw[i]= 1
67          else:
68              yw[i]= -1
69          i+=1
70
71      x=testdataset[k]
72      #calculate x˜ for this sample
73      x_attack= x - P*yw
74      #add to matrix with all x˜
75      x_attack_final[k]=x_attack
76
77  ####################################################
78  #calculate error with attacked input
79  ####################################################
80
81  test_error_rate = evaluate_error_rate(x_attack_final, testlabels, w0_optimal
        , w_optimal)
82
83  print(f"Test dataset error rate with attack vector: {test_error_rate *
        100:.2f}%")
```

**Listing 2:** Python Code to Solve Task 8

## 11.3  Python Code For Task 9

```python
import cvxpy as cp
import numpy as np
import scipy.io
import matplotlib.pyplot as plt


# Open the data file .mat
data = scipy.io.loadmat('classifier_dataset.mat')

# Extract the variables from the data
traindataset = data['traindataset']
testdataset = data['testdataset']


trainlabels = data['trainlabels'].flatten()  # Flatten to convert into 1D
    array
testlabels = data['testlabels'].flatten()    # Flatten to convert into 1D
    array



# Save the number of rows in N (number of samples) and columns in D (number
    of features)
N, D = traindataset.shape

# Regularization parameter
ro = 0.1

# Define the optimization variables
w0 = cp.Variable()          # Bias term (scalar)
w = cp.Variable(D)          # Weights vector (D-dimensional)

# Define the hinge loss function h(u) = max(0, 1 - u)
hinge_losses = cp.pos(1 - cp.multiply(trainlabels, traindataset @ w + w0))

# Define the objective function: hinge loss + regularization term
objective = cp.Minimize((1/N) * cp.sum(hinge_losses) + ro * cp.norm(w, 2)
    **2)

# Define the problem and solve it
problem = cp.Problem(objective)
problem.solve()

# Retrieve the optimal parameters
w_optimal = w.value
w0_optimal = w0.value



# Function to evaluate the classifier error rate fD on a given dataset
def evaluate_error_rate(dataset, labels, w0_opt, w_opt):
```

```
50      predictions = np.sign(dataset @ w_opt + w0_opt)  # Classify
51      misclassifications = np.sum(predictions != labels)  # Count errors
52      error_rate = misclassifications / len(labels)  # Calculate error rate
53      return error_rate
54
55
56  P=0.18
57
58  ##################################################
59  #make the x~ vector (atacker vector)
60  ##################################################
61  N2, D2 = testdataset.shape
62
63  x_attack_final=np.empty(shape=(N2, D2))
64
65
66
67  #For each sample, calculate x~
68  for k in range(0,N2):
69      #calculate y*w
70
71      yw=testlabels[k]*w_optimal
72      i=0
73
74      #aply sign function to all elements in the vector yw
75      for num in yw:
76          if num>=0:
77              yw[i]= 1
78          else:
79              yw[i]= -1
80          i+=1
81
82      x=testdataset[k]
83      #calculate x~ for this sample
84      x_attack= x - P*yw
85      #add to matrix with all x~
86      x_attack_final[k]=x_attack
87
88
89
90  w0 = cp.Variable()
91  w = cp.Variable(D)
92
93  ro = 0.1
94
95
96  #Vector manipulation
97  train_labels_col = cp.reshape(trainlabels, (400,1))
98  w_T = cp.reshape(w, (1,784))
99
100 #Compute yw
101 product_matrix = train_labels_col @ w_T
102
103 #Equivalent to calculating the l1-norm for every row
104 l1_norm = cp.sum(cp.abs(product_matrix), axis=1)
```

```
105
106 #Define the hinge loss function
107 hinge_losses_9 = cp.pos(1 - (cp.multiply(trainlabels, traindataset @ w + w0)
        - P * l1_norm))
108
109 #Problem definition
110 objective_9 = cp.Minimize((1/N) * cp.sum(hinge_losses_9) + ro * cp.norm(w,
        2)**2)
111 problem_9 = cp.Problem(objective_9)
112 problem_9.solve()
113
114 w_optimal = w.value
115 w0_optimal = w0.value
116
117
118 #Result Evaluation
119 train_error_rate = evaluate_error_rate(traindataset, trainlabels, w0_optimal
        , w_optimal)
120 print(f"Training dataset error rate: {train_error_rate * 100:.2f}%")
121
122 test_error_rate = evaluate_error_rate(testdataset, testlabels, w0_optimal,
        w_optimal)
123 print(f"Test dataset error rate: {test_error_rate * 100:.2f}%")
124
125 attack_test_error_rate = evaluate_error_rate(x_attack_final, testlabels,
        w0_optimal, w_optimal)
126 print(f"Test dataset error rate with attack vector: {attack_test_error_rate
        * 100:.2f}%")
```

**Listing 3:** Python Code to Solve Task 9

## 11.4   Python Code For Task 10

```python
 1  import numpy as np
 2  import scipy.io as sio
 3  import matplotlib.pyplot as plt
 4
 5
 6  # Load the .mat file
 7  def load_data(file_path):
 8      data = sio.loadmat(file_path)
 9      X = data['x']  # Input data (X values)
10      Y = data['y']  # Target data (Y values)
11      U = data['u']  # ...
12      V = data['v']  # ...
13      S = data['s']  # ...
14      R = data['r']  # ...
15      return X.flatten(), Y.flatten(), U.flatten(), V.flatten(), S.flatten(),
    R.flatten()
16
17
18
19  # Mixture of Linear Models
20  def mixture_model(x, u, v, s, r, K, N):
21      y_pred = np.zeros(N)
22      W = np.zeros((N, K))
23
24      alphas = np.zeros((N, K))
25      for k in range(K-1):
26          alphas[:, k] = u[k] * x + v[k]
27
28      # Centering to prevent overflow
29      max_alpha = np.max(alphas,axis=1,keepdims=True)
30      exp_alphas = np.exp(alphas - max_alpha)
31
32      # Normalizing weights
33      for k in range(K):
34          W[:, k] = exp_alphas[:, k] / np.sum(exp_alphas, axis=1)
35
36      # Combine predictions from all K models
37      for k in range(K):
38          y_pred += W[:, k] * (s[k] * x + r[k])
39
40      return y_pred, W
41
42
43
44  # Compute partial derivatives (Jacobian elements)
45  def compute_jacobian_partial(x, u, v, s, r, W, k, K, N, param):
46      if param == 's':
47          return W[:, k] * x
48      elif param == 'r':
49          return W[:, k]
50      elif param == 'u':
51          J_partial_u = np.zeros(N)
52          for i in range(N):
```

```
53
54          k_term = W[i, k] * (1 - W[i, k]) * x[i]*(s[k] * x[i] + r[k])#
    derivative of Wk*Yk
55
56          j_term = W[i, k]*x[i]*sum(W[i, j] * (s[j] * x[i] + r[j]) for j
    in range(K) if j != k) # sum of derivative of Wj*Yj
57
58          J_partial_u[i] = k_term - j_term
59
60
61      return J_partial_u
62
63  elif param == 'v':
64      J_partial_v = np.zeros(N)
65      for i in range(N):
66
67          W_term = W[i, k] * (1 - W[i, k])*(s[k] * x[i] + r[k])#
    derivative of Wj*Yj
68
69          sum_term = W[i, k]*sum(W[i, j] * (s[j] * x[i] + r[j]) for j in
    range(K) if j != k) #sum of derivative Wj*Yj
70
71          J_partial_v[i] = W_term - sum_term
72
73      return J_partial_v
74  else:
75      raise ValueError(f"Invalid parameter: {param}")
76
77 # Compute the gradient and Jacobian
78 def compute_Jacobian(x, u, v, s, r, W, K, N):
79     J = np.zeros((N, 4 * K - 2))  # Jacobian
80
81     # Compute partial derivatives for uk and vk
82     for k in range(K - 1):
83         J[:, k] = compute_jacobian_partial(x, u, v, s, r, W, k, K, N, 'u')
    # Jacobian for u_k
84         J[:, (K - 1) + k] = compute_jacobian_partial(x, u, v, s, r, W, k, K,
    N, 'v')  # Jacobian for v_k
85
86     # Compute partial derivatives for sk and rk
87     for k in range(K):
88         J[:, 2*(K - 1)+ k] = compute_jacobian_partial(x, u, v, s, r, W, k, K
    , N, 's')  # Jacobian for s_k
89         J[:, 2*(K - 1)+ K + k] = compute_jacobian_partial(x, u, v, s, r, W,
    k, K, N, 'r')  # Jacobian for r_k
90
91     return J
92
93
94 # Levenberg-Marquardt optimization
95 def levenberg_marquardt(x, y, u, v, s, r, K, N, max_iter=5000, epsilon=1e-4,
    lambda_init=1.0):
96     lambda_ = lambda_init  # Damping parameter
97     residuals_list = []  # To store sum of squared residuals for ploting
98     grad_obj_func_list = [] # To store gradient of objective function for
```

```
      ploting
99
100   for iter in range(max_iter):

101
102       # Compute model prediction and weights
103       y_pred, W = mixture_model(x, u, v, s, r, K, N)
104       # Combine the u, v, s, r vectors into a single column vector
105       param_vector = np.concatenate([u, v, s, r]).reshape(-1, 1)

106
107       # Compute the residual and objective function
108       residual = y_pred - y
109       obj_func = np.sum(residual ** 2)
110       residuals_list.append(obj_func)
111       # Compute the full Jacobian matrix (gradients)
112       J = compute_Jacobian(x, u, v, s, r, W, K, N)

113
114       # Construct the A matrix (Jacobian and regularization term)
115       sqrt_lambda = np.sqrt(lambda_)

116
117       # Create a (4*K-2)x(4*K-2) identity matrix
118       identity_matrix = np.identity(4*K-2)

119
120       # Multiply the identity matrix by sqrt_lambda to get the final
      sqrt_lambda matrix
121       sqrt_lambda_matrix = sqrt_lambda * identity_matrix

122
123       # Combine the jacobian with the square root of the scalar
124       A = np.vstack([J, sqrt_lambda_matrix])

125

126
127       J_param = J @ param_vector  # Multiply jacobian by parameters vector
128       b_top = J_param - residual.reshape(-1, 1)  # Subtract residual from
      each row

129
130       # Bottom part: sqrt_lambda * param_vector
131       b_bottom = sqrt_lambda * identity_matrix @ param_vector  # Shape
      will be 14 x 1

132
133       # Combine the top and bottom parts
134       b = np.vstack([b_top, b_bottom])

135
136       # Solve the least-squares problem for min
137       minimized_parameters,_ , _, _ = np.linalg.lstsq(A, b, rcond=None)

138
139       # Compute the candidate parameters with min
140       u_new, v_new, s_new, r_new = update_parameters(u, v, s, r,
      minimized_parameters, K)
141       # Compute new prediction and objective function with updated
      parameters
142       y_pred_new, W_new = mixture_model(x, u_new, v_new, s_new, r_new, K,
      N)
143       obj_func_new = np.sum((y_pred_new - y) ** 2)

144

145
146       gradient = 2 * J.T @ residual
```

```
147          # Check stopping criterion
148          grad_obj_func_list.append(np.linalg.norm(gradient))
149          if np.abs(np.linalg.norm(gradient))< epsilon:
150              print(f"Converged at iteration {iter}")
151              break
152          # Check if the step is valid
153
154          if obj_func_new < obj_func:  # Valid step
155              u, v, s, r = u_new, v_new, s_new, r_new
156              lambda_ *= 0.7  # Decrease lambda
157
158
159          else:  # Null step
160              lambda_ *= 2.0  # Increase lambda
161
162
163      # Plot the results
164
165      plot_results(x, y, u, v, s, r, y_pred, W,residuals_list,
         grad_obj_func_list)
166      return u, v, s, r
167
168 # Update the parameters (u, v, s, r) with the delta step from LM
169 def update_parameters(u, v, s, r, minimized_parameters, K):
170
171      # Ensure the deltas are treated as 1D vectors instead of 2D
172      u_new = minimized_parameters[:K - 1].flatten()  # Flatten in case it's
         higher dimensional
173      v_new = minimized_parameters[K - 1:2 * K - 2].flatten()
174      s_new  = minimized_parameters[2 * K - 2:3 * K - 2].flatten()
175      r_new  = minimized_parameters[3 * K - 2:].flatten()
176
177      # Update the parameters with the corresponding deltas
178
179
180      return u_new , v_new , s_new , r_new
181
182 # Plot the fitted signal and the weights
183 def plot_results(x, y, u, v, s, r, y_pred, W,residuals_list,
         grad_obj_func_list):
184
185      plt.figure()
186      plt.plot(x, y, 'ko', markersize=3, label='Original Data')
187      plt.plot(x, y_pred, 'o', markerfacecolor='none', markeredgecolor='blue',
         markersize=3, label='Fitted Signal')
188      plt.ylim(min(y) - 1, max(y) + 1)
189      plt.xlim(min(x) - 1, max(x) + 1)
190      plt.title('Fitted Signal using Mixture of Linear Models')
191      plt.legend()
192      plt.grid(True)
193      plt.show()
194
195      plt.figure()
196      for k in range(W.shape[1]):
197          plt.plot(x, W[:, k], label=f'Weight {k+1}')
```

```
198     plt.title('Weights of each Linear Model')
199     plt.legend()
200     plt.grid(True)
201     plt.show()
202
203     plt.figure()
204     plt.plot(range(len(residuals_list)), residuals_list, 'b-', label='
        Objective funtion')
205     plt.xlabel('Iterations')
206     plt.ylabel('Objective funtion')
207     plt.title('Objective funtion Over Iterations')
208     plt.legend()
209     plt.grid(True)
210     plt.yscale('log')
211     plt.show()
212
213     plt.figure()
214     plt.plot(range(len(grad_obj_func_list)), grad_obj_func_list, 'b-', label
        ='Gradient of objective funtion')
215     plt.xlabel('Iterations')
216     plt.ylabel('Gradient of objective funtion')
217     plt.title('Gradient of objective funtion Over Iterations')
218     plt.legend()
219     plt.ylim(0.00001, 100000)
220     plt.grid(True)
221     plt.yscale('log')
222     plt.show()
223
224
225
226
227 # Example usage
228 file_path = 'lm_dataset_task.mat'  # Replace with actual path
229 # Initial values for u, v, s, r
230 x, y, u, v, s, r = load_data(file_path)
231
232
233 K = len(s)  # Number of models
234 N = len(x)  # Number of data points
235
236 # Run the LM optimization
237 u_opt, v_opt, s_opt, r_opt = levenberg_marquardt(x, y, u, v, s, r, K, N)
238 #print results
239 print(u_opt)
240 print(v_opt)
241 print(s_opt)
242 print(r_opt)
```

**Listing 4:** Python Code to Solve Task 10